# The AI black box explanation problem

## Dino Pedreschi

## KDD Lab Dip. Informatica, Univ. Pisa & ISTI-CNR

**www.sobigdata.eu**

**DATA**

**ANALYTICS**

**COMPUTING INFRASTRUCTURE**

**DATA ETHICS**

**DATA SCIENCE**

**SCIENCE**
- Astronomy
- Social sciences
- Medicine
- Environment
- Agrifood
- ...

**SOCIETY**
- Policy
- Citizens
- Social Good
- ...

**INDUSTRY & BUSINESS**
- Media, Entertainment and Information,
- Consumer,
- Healthcare,
- Energy
- Information and Communication Techno...
- Mobility
- Financial Services
- ...

, Giannotti, Pedreschi, G7 2017 Position paper on Data Science
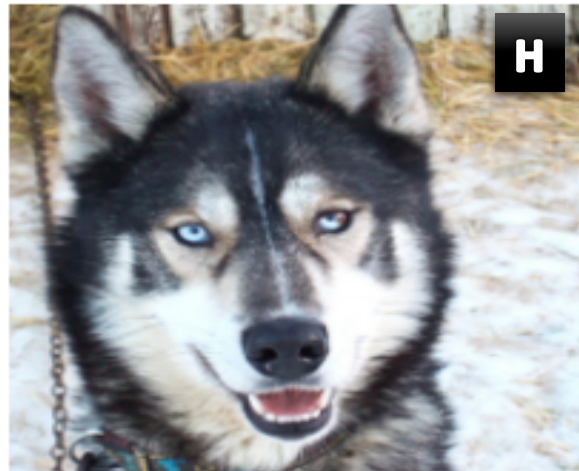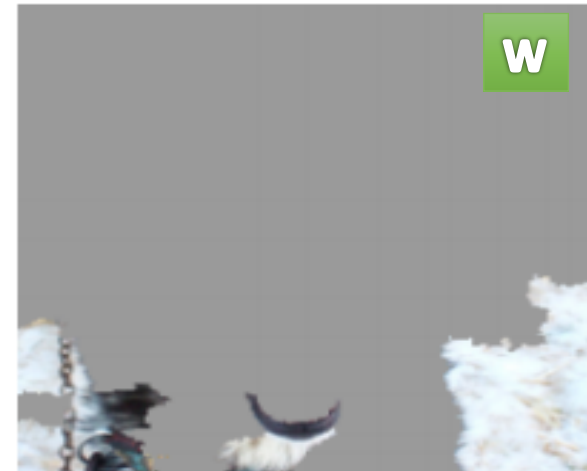
Deep learning classifies a wolf in a picture based on ... **the snow in the background!**



(a) Husky classified as wolf  (b) Explanation

Ribeiro et al., KDD 2016

ilitary tank classification depends on the backgroun

# COMPAS crime recidivism score



**DYLAN FUGETT**

Prior Offense
1 attempted burglary

Subsequent Offenses
3 drug possessions

LOW RISK    3

**BERNARD PARKER**

Prior Offense
1 resisting arrest
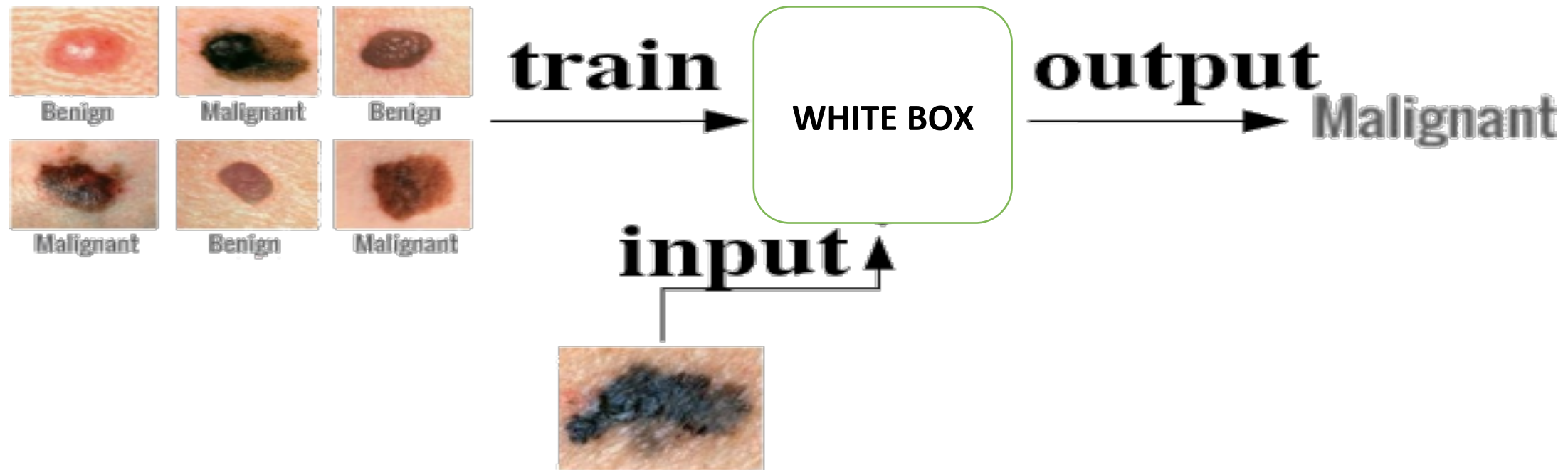without violence
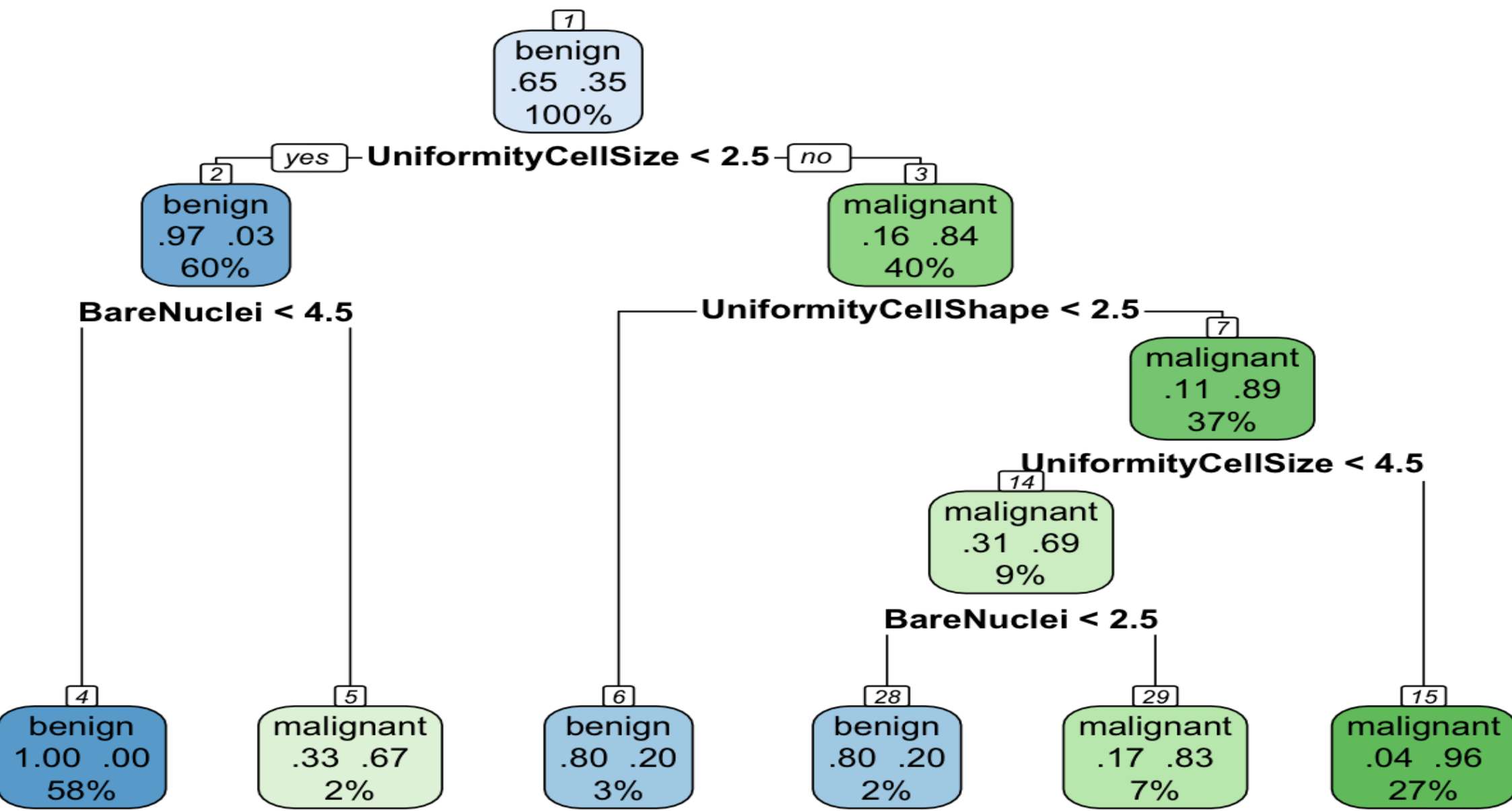
Subsequent Offenses
None

HIGH RISK    10

*Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.*
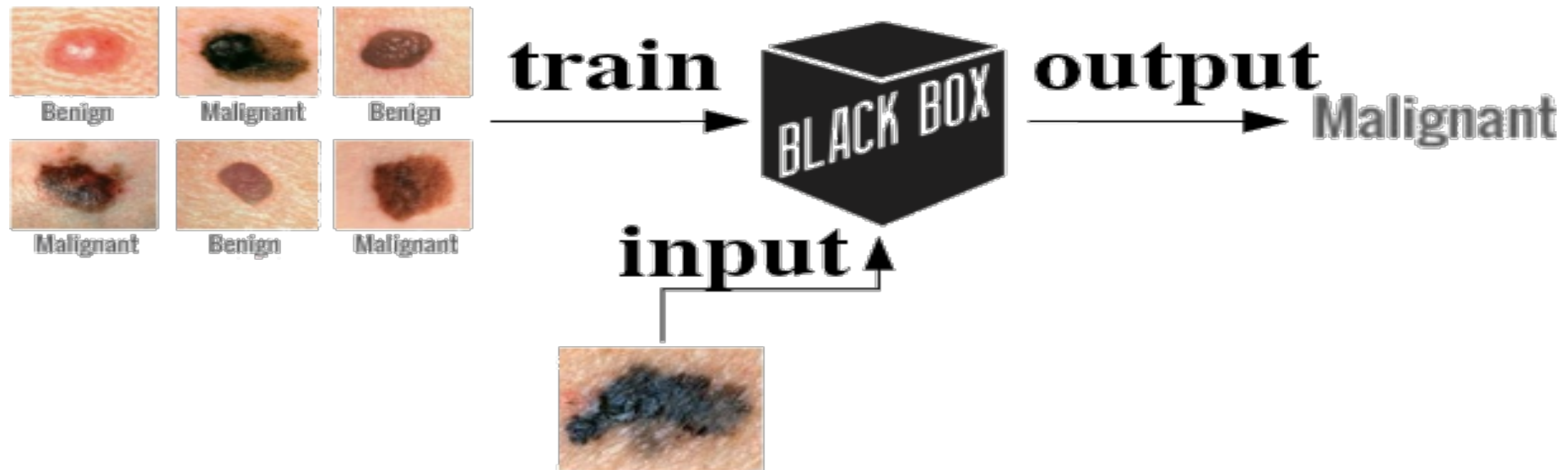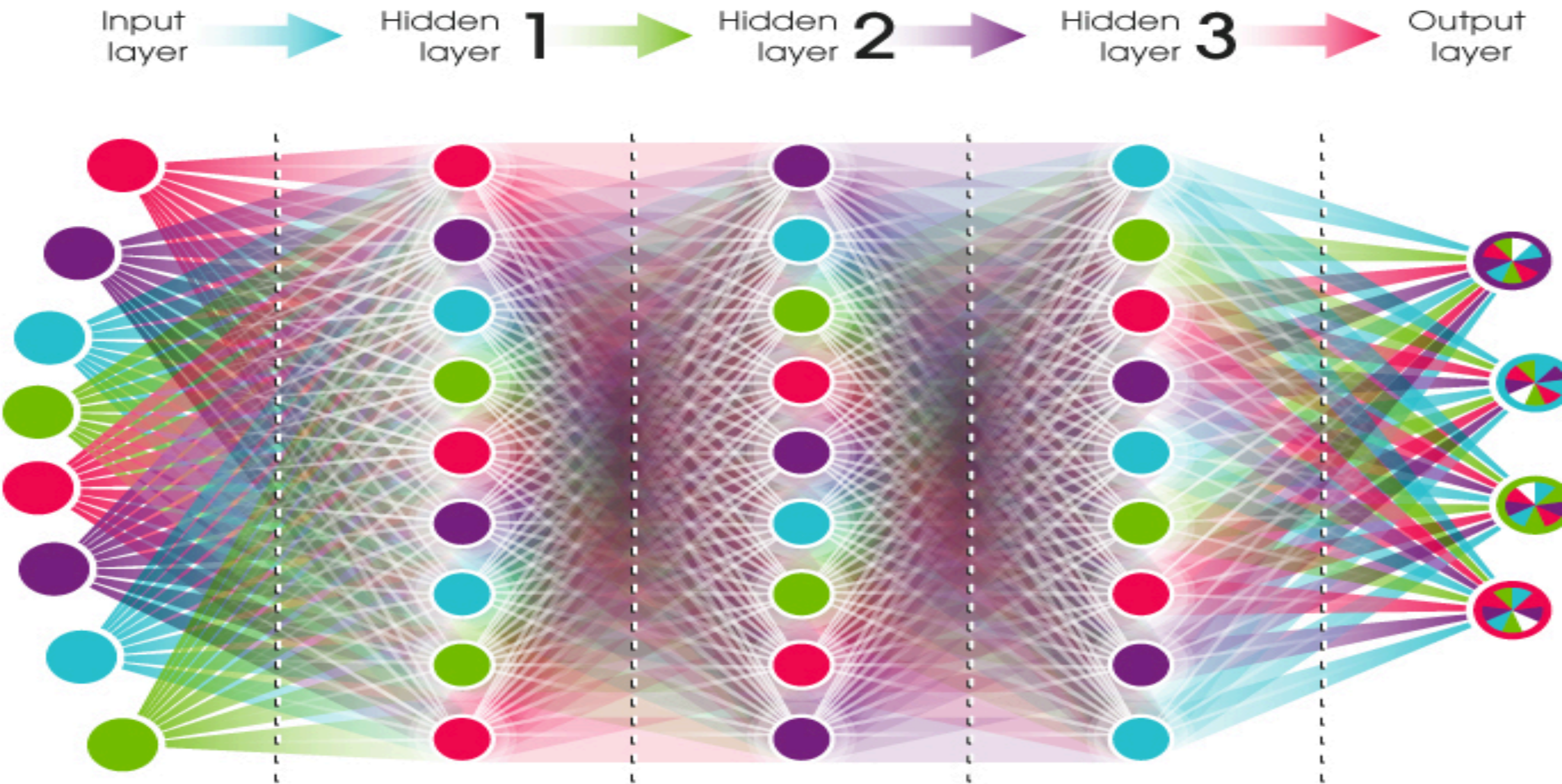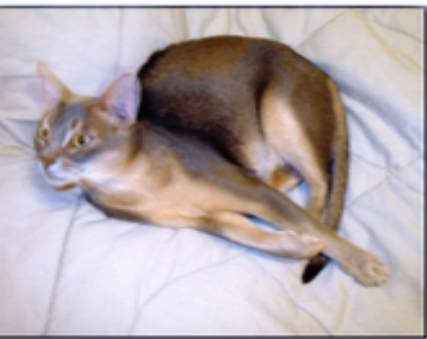
No Amazon free same-day delivery
for restricted minority neighborhoods

Black res...

White residents    Black residents

Same-day
delivery
area

...sidents    Black res...

Same-day
delivery
area

Same-day
delivery
area

...dents    Black res...

Same-day
delivery area

Milwaukee
1,342,594

Indianapolis
1,566,866

Chicago
4,228,361

Cincinn...
1,390...

Louisville
910,289

Nashville
1,316,372

New Y...
to Phila...
14.73...

Atla...
1,601...

Fresno area
1,374,505

...o area
...,015

...Angeles area
...,968,496

San Diego
2,317,377

Phoenix
2,886,340

Tucson
922.273

Dallas &
Fort Worth
3,570,116

Tampa Bay ar...
1,671,604

loomberg analyis of data from Amazon.com
merican Community Survey

# AI = Machine Learning + Big Data

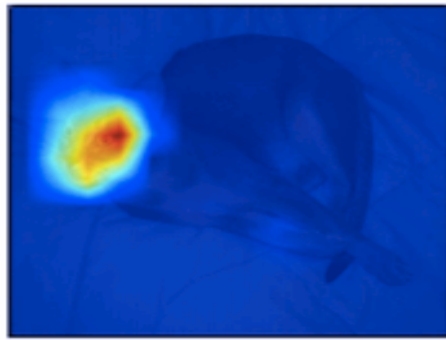# AI = Machine Learning + Big Data

# DEEP NEURAL NETWORK

Input layer → Hidden layer 1 → Hidden layer 2 → Hidden layer 3 → Output layer

# We risk to create algorithms we don't fully understand

**...stems that recommend humans making a decision... ...ould explain why**

...requires **transparency** and **responsibility**

...E.g., in healthcare, justice, surveillance, predictive ...policing, autonomous vehicles, risk scoring, ...

ature | International weekly journal of science

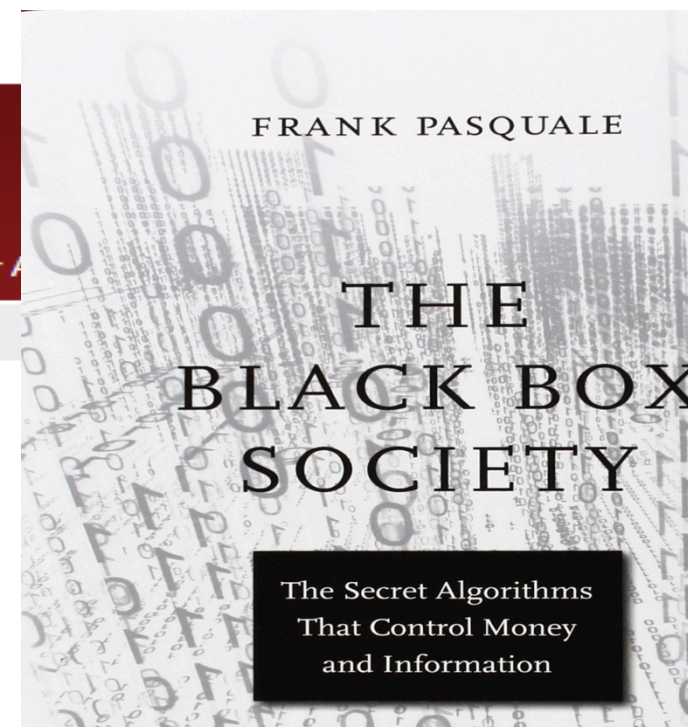News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video | For A...

Volume 537 | Issue 7621 | Editorial | Article

E | EDITORIAL

...e accountability for big-data algorithms

...d bias and improve transparency, algorithm designers must make data sources and ...s public.

...ember 2016

FRANK PASQUALE

THE BLACK BOX SOCIETY

The Secret Algorithms That Control Money and Information

# ght of explanation

Since 25 May 2018, GDPR establishes a right for all individuals to obtain "**meaningful explanations of the logic involved**" when "**automated (algorithmic) individual decision making**", including profiling, takes place.